

# Revisiting maximum-a-posteriori estimation in log-concave models: from differential geometry to Bayesian decision theory

M. Pereyra \*

January 24, 2017

## Abstract

Maximum-a-posteriori (MAP) estimation has become the main Bayesian estimation methodology in many areas of modern data science such as mathematical imaging and machine learning, where high dimensionality is addressed by using models that are log-concave and whose posterior mode can be computed very efficiently by using convex optimisation algorithms. However, despite its success and rapid adoption, MAP estimation is not theoretically well understood yet, and the prevalent view is that it is generally not proper Bayesian estimation in a decision-theoretic sense. This paper presents a new decision-theoretic derivation of MAP estimation in Bayesian models that are log-concave. Our analysis is based on differential geometry and proceeds as follows. First, we exploit the log-concavity of the model to induce a Riemannian geometry on the parameter space. We then use differential geometry to identify the natural or canonical loss function to perform Bayesian point estimation in that Riemannian manifold. For log-concave models this canonical loss is the Bregman divergence of the negative log posterior density, a similarity measure rooted in convex analysis that in addition to the relative position of points also takes into account the geometry of the space, and which generalises the Euclidean squared distance to non-Euclidean settings. We then show that the MAP estimator is the Bayesian estimator that minimises the expected canonical loss, and that the posterior mean or minimum mean squared error (MMSE) estimator is the Bayesian estimator that minimises the dual canonical loss. Finally, we establish universal performance and stability guarantees for MAP and MMSE estimation in high dimensional log-concave models. These theoretical results provide a new understanding of MAP and MMSE estimation under log-concavity, and reveal new insights about their good empirical performance and about the roles that log-concavity plays in high dimensional inference problems.

---

\*School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh EH14 4AS, United Kingdom. Email: m.pereyra@hw.ac.uk.

# 1 Introduction

We consider the Bayesian estimation of an unknown quantify of interest  $x \in \mathbb{R}^n$  from an observation  $y$  [18]. We focus on Bayesian models whose posterior distribution is log-concave, i.e.,

$$p(x|y) = \frac{\exp\{-\phi(x)\}}{\int_{\mathbb{R}^n} \exp\{-\phi(s)\} ds}, \quad (1)$$

for some proper convex function  $\phi : \mathbb{R}^n \rightarrow (-\infty, \infty]$ . Such models are ubiquitous in modern data science, particularly for problems where  $x$  is high dimensional (e.g.,  $n > 10^3$ ). For example, statistical imaging and machine learning methods rely strongly on log-concave models of the form  $\phi(x) = \|y - Ax\|^2/2\sigma^2 + \phi(Bx) + \mathbf{1}_{\mathcal{S}}(x)$  for some linear operators  $A$  and  $B$ , convex regulariser  $\phi$ , and convex set constraint  $\mathcal{S}$ , and are typically of dimension  $n > 10^5$  [17, 14, 11].

Because drawing conclusions directly from  $p(x|y)$  is difficult, Bayesian methods generally deliver summaries of  $p(x|y)$ , namely Bayes point estimators, which summarises the information in  $p(x|y)$  optimally in the following decision-theoretic sense [18]:

**Definition 1.1.** Let  $L : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  be a loss function that quantifies the difference between two points in  $\mathbb{R}^n$ . A Bayes estimator associated with  $L$  is any estimator that minimises the posterior expected loss, i.e.,

$$\hat{x}_L = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[L(u, x)] \triangleq \int_{\mathbb{R}^n} L(u, x) p(x|y) dx.$$

The loss function  $L$  is specified by the analyst and usually verifies the following general conditions:

- $L(u, x) \geq 0, \forall u, x \in \mathbb{R}^n$ ,
- $L(u, x) = 0 \iff u = x$ ,
- $L$  strictly convex w.r.t. its first argument (to guarantee estimator uniqueness).

Observe that  $L$  is not necessarily symmetric, i.e.,  $L(u, x) \neq L(x, u)$ . We do not enforce symmetry because the arguments of  $L$  have clearly different roles in the decision problem.

Ideally  $L$  should be chosen carefully based on specific aspects of the problem and application considered. This is particularly important for instance in imaging problems that are severely ill-posed or ill-conditioned, where this choice can significantly impact estimation results. However, specifying a bespoke loss function for high dimensional problems is not easy, and as a result most methods reported in the literature use default losses and estimators.

In particular, Bayesian methods in engineering fields such as imaging have traditionally used the minimum mean squared error (MMSE) estimator, which is

given by the posterior mean  $\hat{x}_{MMSE} = \int_{\mathbb{R}^n} p(x|y) x dx$ . This estimator is widely regarded as a gold standard in these fields, in part because of its good empirical performance and favourable theoretical properties, and also perhaps in part because of cultural heritage. From Bayesian decision theory, MMSE estimation is optimal with respect to the entire class of quadratic loss functions of the form  $L(u, x) = (u - x)^\top Q(u - x)$  with  $Q \in \mathbb{R}^{n \times n}$  positive definite [18]. This class provides a second order approximation to any strongly convex loss function, and hence  $\hat{x}_{MMSE}$  is also a proxy for other Bayesian estimators. Also, the quadratic loss is directly related to the Euclidean squared distance, giving  $\hat{x}_{MMSE}$  a clear geometric interpretation. In addition, it has been established in [3] that  $\hat{x}_{MMSE}$  is also optimal w.r.t. the second argument of any Bregman divergence (i.e., any loss function of the form  $D_h(x, u) = h(x) - h(u) - \nabla h(u)^\top (x - u)$  for a convex function  $h \in \mathcal{C}^1$ ), a more general class of loss functions that includes quadratic losses and that is related to non-Euclidean geometries [1].

Unfortunately, calculating  $\hat{x}_{MMSE}$  in high dimensional models can be very difficult because it requires solving integrals that are often too computationally expensive for the applications considered. This has stimulated much research on the topic, from fast Monte Carlo simulation methods to efficient approximations with deterministic algorithms [17]. But with ever increasingly large problems and datasets, many applied fields have progressively focused on alternatives to MMSE estimation.

In particular, modern imaging and machine learning methods rely strongly on maximum-a-posteriori (MAP) estimation

$$\begin{aligned}\hat{x}_{MAP} &= \operatorname{argmax}_{x \in \mathbb{R}^n} p(x|y), \\ &= \operatorname{argmin}_{x \in \mathbb{R}^n} \phi(x),\end{aligned}$$

whose calculation is a convex problem that can often be solved very efficiently, even in very high dimensions (e.g.,  $n > 10^7$ ), by using convex optimisation techniques [7, 12, 14]. Interestingly, modern non-statistical imaging and machine learning methods also predominately solve problems by convex optimisation, and their solutions are often equivalent to performing MAP estimation for some implicit Bayesian model.

There is abundant evidence that MAP estimation delivers accurate results for a wide range of log-concave models and applications. However, from a theoretical viewpoint MAP estimation is not well understood. Currently the predominant view is that MAP estimation is not formal Bayesian estimation in the decision-theoretic sense postulated by Definition 1.1 because it does not minimise a known expected loss. The prevailing interpretation is that MAP estimation is in fact an approximation arising from the degenerate loss  $L_\epsilon(u, x) = \mathbf{1}_{\|x-u\| < \epsilon}$  with  $\epsilon \rightarrow 0$  [18] (this derivation holds for all log-concave models, but is not generally true [4]). However, this asymptotic derivation does not lead to a proper Bayesian estimator. More importantly, the resulting loss is very difficult to motivate for inference problems in  $\mathbb{R}^n$ , and does not explain the good empirical performance reported in the literature.

Furthermore, most other theoretical results for MAP estimation only hold for very specific models, or have been derived by adopting analyses that are extrinsic to the Bayesian decision theory framework (e.g. by analysing MAP estimation as constrained or regularised least-squares regression [9, 10]). As a trivial example, when  $p(x|y)$  is symmetric we have  $\hat{x}_{MAP} = \hat{x}_{MMSE}$ , and thus MAP estimation inherits the favourable properties of MMSE estimation. This result has been partially extended to some denoising models of the form  $p(x|y) \propto \exp\{\|y - x\|^2/2\sigma^2 + \lambda h(x)\}$  in [15], where it is shown that MAP estimation coincides with MMSE estimation with a different model  $\tilde{p}(x|y) \propto \exp\{\|y - x\|^2/2\sigma^2 + \tilde{\lambda}\tilde{h}(x)\}$ . It follows that for these models MAP estimation is decision-theoretic Bayesian estimation w.r.t. the weighted loss  $L(u, x) = \|u - x\| \exp\{\tilde{\lambda}\tilde{h}(x) - \lambda h(x)\}$ . This is a post-hoc loss, but the result is interesting because it highlights that a single estimator may have a plurality of origins. Lastly, Burger & Lucka [8] recently established that MAP estimation is decision-theoretic Bayesian estimation for all linear Gaussian models of the form  $p(x|y) \propto \exp\{\|y - Ax\|_{\Sigma^{-1}}^2/2 + \lambda h(x)\}$ , where  $A$  is a known linear operator,  $\Sigma$  a known noise covariance, and  $h$  is convex and Lipschitz continuous. More precisely, that paper shows that for these models MAP estimation is optimal w.r.t. the loss  $L(u, x) = \|A(u - x)\|_{\Sigma^{-1}}^2 + 2\lambda D_h(u, x)\}$ , where  $D_h(x) = h(u) - h(x) - \nabla h(x)^\top(u - x)$  is the Bregman divergence associated with  $h$ . It may appear that this loss is rather artificial and difficult to analyse and motivate, however the new results presented in section 3 show that it is a specific instance of a more general loss that stems directly from the consideration of the model geometry.

In order to understand MAP estimation, in this paper we first revisit the choice of the loss function for Bayesian point estimation in the context of log-concave models. A main novelty is that, instead of specifying the loss directly, we use differential geometry to derive it automatically from the geometry of the model. Precisely, we show that under some regularity assumptions, the log-concavity of  $p(x|y)$  induces a specific Riemannian differential geometry on the parameter space, and that taking into account this space geometry naturally leads to an intrinsic or canonical loss function to perform Bayesian point estimation. Following on from this, we establish that the canonical loss for the parameter space is the Bregman divergence associated with  $\log p(x|y)$ , and that the Bayesian estimator w.r.t. this loss is the MAP estimator. We then show that the MMSE estimator is the Bayesian estimator associated with the dual canonical loss, and propose universal estimation performance guarantees for MAP and MMSE estimation in log-concave models.

The remainder of the paper is organised as follows: section 2 introduces the elements of differential geometry that are essential to our analysis. In section 3 we present our main theoretical result: a decision-theoretic and differential-geometric derivation of MAP and MMSE estimation, as well as universal bounds on the estimation error involved. Conclusions are finally reported in section 4. Proofs are presented in the appendix.

## 2 Riemannian geometry and the canonical divergence

In this section we recall some elements of differential geometry that are necessary for our analysis. For a detailed introduction to this topic we refer the reader to [1].

An  $n$ -dimensional Riemannian manifold  $(\mathbb{R}^n, g)$ , with metric  $g : \mathbb{R}^n \rightarrow \mathcal{S}_{++}^n$  and global coordinate system  $x$ , is a vector space that behaves locally as an Euclidean space. Precisely, for any point  $x \in \mathbb{R}^n$  we have a tangent space  $\mathcal{T}_x \mathbb{R}^n$  with inner product  $\langle u, x \rangle = u^\top g(x)x$  and norm  $\|x\| = \sqrt{x^\top g(x)x}$ . This geometry is local and may vary smoothly from  $\mathcal{T}_x M$  to neighbouring tangent spaces. These variations are encoded in the manifold's affine connection  $\Gamma$ , with coefficients given by  $\Gamma_{ij,k}(x) = \partial_k g_{i,j}(x)$ .

Moreover, similarly to Euclidean spaces, the manifold  $(\mathbb{R}^n, g)$  supports divergence functions.

**Definition 2.1.** A function  $D : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  is a divergence function on  $\mathbb{R}^n$  if the following conditions hold for any  $u, x \in \mathbb{R}^n$ :

- $D(u, x) \geq 0, \forall u, x \in \mathbb{R}^n$ ,
- $D(u, x) = 0 \iff x = u$ ,
- $D(u, x)$  is strongly convex w.r.t.  $u$ , and  $\mathcal{C}^2$  w.r.t  $u$  and  $x$ .

The class of divergence functions coincides with that of loss functions for Bayesian point estimation considered in section 1 with mild additional regularity conditions. Hence divergence functions are sensible candidates to define Bayesian estimators. Divergence functions also provide a link to the differential geometry of the space, which allows relating space geometry and Bayesian decision theory. This relationship has been used previously to analyse Bayesian decision problems from a Riemannian geometric viewpoint, leading to the so-called decision geometry framework [13]. Here we adopt an opposite perspective: we start by considering a Riemannian manifold  $(\mathbb{R}^n, g)$  and then use the relationship to identify the divergence functions that arise naturally in that space. In particular, we focus on the so-called canonical divergence on  $(\mathbb{R}^n, g)$ , which is a generalisation of the Euclidean squared distance to this kind of manifold [2].

**Definition 2.2** (Canonical divergence [2]). For any two points  $u, x \in \mathbb{R}^n$ , the  $(\mathbb{R}^n, g)$ -canonical divergence is given by

$$D(u, x) = \int_0^1 \int_0^1 t \dot{\gamma}_t^\top g(\gamma_t) \dot{\gamma}_t dt, \quad (2)$$

where  $\gamma_t$  is the  $\Gamma$ -geodesic connecting  $u \rightarrow x$  and  $\dot{\gamma}_t = d/dt \gamma_t$ .

To gain a geometric intuition for  $D$  it is useful to compare it to the length of the  $\Gamma$ -geodesic between  $u$  and  $x$ . Precisely, by noting that the squared length of a curve  $\zeta_t : [0, 1] \rightarrow \mathbb{R}^n$  on the manifold  $(\mathbb{R}^n, g)$  is given by  $\int_0^1 \dot{\zeta}_t^\top g(\zeta_t) \dot{\zeta}_t dt$ , we observe

that  $D(u, x)$  is essentially the squared length of the  $\Gamma$ -geodesic  $\gamma_t$  weighted linearly along the path from  $u$  to  $x$ . This weighting in (2) guarantees that  $D(u, x)$  is convex in  $u$ , a necessary condition to define a divergence function (the weighting also leads to other important properties such as linearity w.r.t.  $g$ , see section 3). The linear weighting also introduces an asymmetry, i.e., generally  $D(u, x) \neq D(x, u)$ , which will have deep implications for Bayesian estimation.

Finally, it is easy to check that (2) reduces to the Euclidean squared distance  $D(u, x) = \frac{1}{2}(u - x)^\top g(u - x)$  when  $(\mathbb{R}^n, g)$  is the Euclidean space with inner product  $\langle u, x \rangle = u^\top g x$ . More generally,  $D$  is always consistent with the local Euclidean geometry of the manifold  $(\mathbb{R}^n, g)$ . That is, for any point  $x + dx$  in the neighbourhood of  $x$  we have  $D(x + dx, x) = \|dx\|^2/2 + o(\|dx\|^2)$ , where  $\|\cdot\|$  is the Euclidean norm of the tangent space  $\mathcal{T}_x \mathbb{R}^n$  (a higher order approximation of  $D(x + dx, x)$  is also possible by using the affine connection  $\Gamma$  [1]). And if we use the decision geometry framework [13] to derive the Riemannian geometry induced by  $D$  on  $\mathbb{R}^n$  we obtain

$$g_{i,j}^{(D)}(x) \triangleq \partial_i \partial_j D(x, x) = g_{i,j}(x), \quad \Gamma_{ij,k}^{(D)}(x) \triangleq \partial_i \partial_j \partial'_k D(x, x) = \Gamma_{ij,k}(x),$$

(here  $\partial$  and  $\partial'$  denote differentiation w.r.t. the first and second components of  $D$  respective), which indicates that  $D$  is fully specified by  $(\mathbb{R}^n, g)$  and that it induces the same space geometry that originated it in the first place.

### 3 A geometric derivation of MAP and MMSE estimation

#### 3.1 From differential geometry to Bayesian decision theory

In this section we use differential geometry to relate  $p(x|y)$  to the loss functions used for Bayesian point estimation of  $x$ . Precisely, we exploit the log-concavity of  $p(x|y)$  to induce a Riemannian geometry on the parameter space. This in turn defines a canonical loss for that space and two Bayesian estimators w.r.t. to this loss: a primal estimator related to  $D(u, x)$  and a dual estimator related to the dual divergence  $D_\phi^*(u, x) = D_\phi(x, u)$ . We first consider the case where  $p(x|y)$  is smooth and strongly log-concave, and later analyse the effect of relaxing these assumptions.

**Theorem 3.1** (Canonical Bayesian estimators). *Suppose that  $\phi(x) = -\log \pi(x|y)$  is strongly convex, continuous, and  $\mathcal{C}^3$  on  $\mathbb{R}^n$ . Let  $(\mathbb{R}^n, g)$  denote the Riemannian manifold induced by  $\phi$ , with metric coefficients  $g_{i,j}(x) = \partial_i \partial_j \phi(x)$ . Then, the canonical divergence on  $(\mathbb{R}^n, g)$  is the Bregman divergence associated with  $\phi$ , i.e.,*

$$D_\phi(u, x) = \phi(u) - \phi(x) - \nabla \phi(x)(u - x).$$

*In addition, the Bayesian estimator associated with  $D_\phi(u, x)$  is unique and is given*

by the maximum-a-posteriori estimator,

$$\begin{aligned}\hat{x}_{D_\phi} &\triangleq \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(u, x)], \\ &= \operatorname{argmin}_{x \in \mathbb{R}^n} \phi(x), \\ &= \hat{x}_{MAP}.\end{aligned}$$

The Bayesian estimator associated with the dual canonical divergence  $D_\phi^*(u, x) = D_\phi(x, u)$  is also unique and is given by the minimum mean squared error estimator

$$\begin{aligned}\hat{x}_{D_\phi^*} &\triangleq \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi^*(u, x)], \\ &= \int_{\mathbb{R}^n} x p(x|y) dx, \\ &= \hat{x}_{MMSE}.\end{aligned}$$

The proof is reported in the appendix.

The way in which the Bregman divergence  $D_\phi(u, x)$  measures the similarity between  $u$  and  $x$  is directly related to the log-concavity of  $p(x|y)$ . Precisely, because  $\phi$  is strongly convex, then  $\phi(u) > \phi(x) - \nabla \phi(x)(u - x)$  for any  $u \neq x$ . The divergence  $D_\phi(u, x)$  essentially quantifies this gap, which as mentioned previously is directly related to the length of the affine geodesic from  $u$  to  $x$  (and hence not only to the relative position of  $u$  and  $x$  but also to the space geometry). Moreover,  $D_\phi$  is linear w.r.t.  $\phi$ . For example, if  $\phi = \alpha\phi_1 + \beta\phi_2$  for two convex functions  $\phi_1$  and  $\phi_2$  and  $\alpha, \beta \in \mathbb{R}$ , then  $D_\phi = \alpha D_{\phi_1} + \beta D_{\phi_2}$ ; it follows that for the specific case of Gaussian observation models  $D_\phi$  is equivalent to the loss identified in [8].

Finally, Theorem 3.1 provides several valuable new insights into MAP and MMSE estimation in log-concave models. First, MAP estimation stems from Bayesian decision theory, and hence it stands on the same theoretical footing as the core Bayesian methodologies such as MMSE estimation. Second, the conventional definition of the MAP estimator as the maximiser  $\hat{x}_{MAP} = \operatorname{argmax}_{x \in \mathbb{R}^n} p(x|y)$  is mainly algorithmic for these models, useful to highlight that these estimators take the form of a convex optimisation problem that can be solved efficiently by convex optimisation (which is an important computational advantage over other Bayesian point estimators). Third, Theorem 3.1 also reveals a surprising form of duality between MAP and MMSE estimation, which are intimately related to each other by the (asymmetry of the) canonical divergence function that  $p(x|y)$  induces on the parameter space. This duality also manifests itself in other ways. For example it is easy to show that  $\hat{x}_{MMSE}$  is the Bayesian estimator associated with  $D_{\phi^*}(u, x)$ , where  $\phi^*(s) = \max_{x \in \mathbb{R}^n} x^\top s - \phi(x)$  is the convex dual or convex conjugate of  $\phi$  (see the appendix for details). Similarly, noting that  $\partial_i \partial_j \phi^*(x) = g_{i,j}^{-1}(x)$  we see that  $\hat{x}_{MMSE}$  plays the role of  $\hat{x}_{MAP}$  on the manifold  $(\mathbb{R}^n, g^{-1})$ . The case of Gaussian models is particular because  $(\mathbb{R}^n, g)$  is Euclidean, which is a self-dual space; consequently  $D_\phi(u, x) = D_\phi(x, u) = \frac{1}{2} \|u - x\|_{\Sigma^{-1}}^2$  and the primal and

dual canonical estimators coincide. Finally, Theorem 3.1 also shows that under log-concavity and smoothness the posterior mode is a global property of  $p(x|y)$ , which is otherwise not an intuitive property.

### 3.2 Error bounds for MAP and MMSE estimation

We now establish performance guarantees for MAP and MMSE estimation when  $p(x|y)$  is log-concave. Precisely, we establish universal estimation error bounds w.r.t. the dual error function  $D_\phi^*(s, x)$ . Here we do not assume that  $\phi$  is smooth; if  $\phi \notin \mathcal{C}^1$  we replace  $D_\phi^*(s, x)$  with the generalised divergence  $D_{\phi, q}^*(s, x) = \phi(x) - \phi(s) - q^\top(x - s)$  where  $q \in \mathbb{R}^n$  belongs to the subdifferential set of  $\phi$  at  $s$  [5]. We first present the following universal bounds on the expected estimation error:

**Proposition 3.1** (Expected error bound). *Suppose that  $\phi(x) = -\log \pi(x|y)$  is convex on  $\mathbb{R}^n$ . Then,*

$$\mathbb{E}_{x|y} \left[ \frac{D_{\phi, 0}^*(\hat{x}_{MAP}, x)}{n} \right] \leq 1.$$

*In addition, if  $\phi \in \mathcal{C}^1$  then*

$$\mathbb{E}_{x|y} \left[ \frac{D_\phi^*(\hat{x}_{MMSE}, x)}{n} \right] \leq \mathbb{E}_{x|y} \left[ \frac{D_\phi^*(\hat{x}_{MAP}, x)}{n} \right] \leq 1.$$

*Proof.* The proof is reported in the appendix.

Theorem 3.1 establishes that  $\hat{x}_{MMSE}$  minimises the expected dual canonical loss  $D_\phi^*$ , and Proposition 3.1 complements this result by providing an explicit and general upper bound on the loss incurred by using this Bayesian estimator. Proposition 3.1 also states that this bound also applies to  $\hat{x}_{MAP}$ , and that the expected loss per coordinate (e.g., per pixel in imaging problems) cannot exceed 1. Observe that this is also a high dimensional stability result for MAP and MMSE estimation, which provides a theoretical argument for their good empirical performance in imaging, machine learning, and other large scale problems.

Moreover, we also have the following universal large error bound for MAP estimation:

**Proposition 3.2** (Large error bound). *Suppose that  $\phi(x) = -\log \pi(x|y)$  is convex on  $\mathbb{R}^n$ . Then, for any  $\epsilon \in (0, \frac{4}{\sqrt{3}})$*

$$\mathbb{P} \left[ \frac{D_{\phi, 0}^*(\hat{x}_{MAP}, x)}{n} \geq 1 + \epsilon \mid y \right] \leq 3e^{-\frac{n\epsilon^2}{16}}.$$

*Proof.* The proof is reported in the appendix.

Proposition 3.2 essentially indicates that in high dimensional settings the true value of  $x$  is almost certainly within the set  $\{x : D_\phi^*(\hat{x}_{MAP}, x)n^{-1} < 1\}$ , because the probability of a larger error vanishes exponentially fast as  $n$  increases.



Again, this theoretical result supports the vast empirical evidence that MAP estimation delivers accurate results in large-scale convex problems. It also follows from Proposition 3.2 that in such problems  $\hat{x}_{MAP}$  and  $\hat{x}_{MMSE}$  are close to each other (i.e., that  $D_\phi^*(\hat{x}_{MAP}, \hat{x}_{MMSE})n^{-1} \leq 1$  with high probability).

### 3.3 Relaxation of regularity conditions

We now examine the effect of relaxing the regularity assumptions of Theorem 3.1. We consider three main cases: lack of smoothness, lack of strong convexity, and lack of continuity.

#### 3.3.1 Non-smooth models

Several models used in imaging and machine learning are not smooth because they involve priors with non-differentiable points, such as priors based on the  $\ell_1$  norm, the nuclear norm, and the total-variation pseudo-norm [14]. The results of Theorem 3.1 hold for these models with the following minor modifications.

First, observe that these non-smooth models are  $\mathcal{C}^3$  almost everywhere; that is, the set of non-differentiable points has dimension  $n - 1$ , and consequently it has no probability mass and can be omitted in the computation of expectations. Second, because the non-differentiable points do not have Euclidean tangent spaces, instead of a global manifold we consider the collection local manifolds associated with the regions of  $\mathbb{R}^n$  where  $p(x|y)$  is  $\mathcal{C}^3$ . Each one of these regions has a local canonical divergence given by the Bregman divergence  $D(u, x) = D_\phi(u, x) = \phi(u) - \phi(x) - \nabla\phi(x)^\top(u - x)$ . Therefore, for these models we posit  $D_\phi(u, x)$  as the global loss function for any  $(u, x) \in \mathbb{R}^n \times \mathbb{R}^n$  [technically the global loss is the generalised Bregman divergence  $D_\phi(u, x) = \phi(u) - \phi(x) - q_x^\top(u - x)$ , where  $q_x$  belongs to the subdifferential set of  $\phi$  at  $x$  [5], however the expectation  $\mathbb{E}_{x|y}[D_\phi(u, x)]$  is taken over the points where  $\phi$  is  $\mathcal{C}^3$  and hence  $q_x = \nabla\phi(x)$ ]. We then consider the primal and dual Bayesian estimators related to this global loss and obtain that  $\hat{x}_{MAP} = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(u, x)]$  and  $\hat{x}_{MMSE} = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi^*(u, x)]$ , similarly to Theorem 3.1. Observe that in the computation of  $\hat{x}_{MAP}$  and  $\hat{x}_{MMSE}$  the argument  $u$  is optimised over  $\mathbb{R}^n$  including non-differentiable points.

Finally, we note that despite not being a global canonical divergence,  $D_\phi(u, x)$  is still consistent with the space's Riemannian geometry which is local. In addition, the key high dimensional performance guarantees of Propositions 3.1 and 3.2 also hold because  $\phi$  is convex.

#### 3.3.2 Strictly log-concave models

For models that are strictly log-concave but not strongly log-concave only the second and third results of Theorem 3.1 remain true. It is easy to check that the Bayesian estimator associated with  $D_\phi$  is  $\hat{x}_{MAP} = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(u, x)]$ , and that  $\hat{x}_{MMSE} = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi^*(u, x)]$ , similarly to strongly log-concave

models. Therefore, the decision-theoretic derivation of  $\hat{x}_{MAP}$  remains valid, and  $\hat{x}_{MAP}$  and  $\hat{x}_{MMSE}$  remain dual to each other. The high dimensional performance guarantees of Propositions 3.1 and 3.2 also hold because  $\phi$  is convex. However, without strong convexity,  $g$  becomes semi-positive definite and  $(\mathbb{R}^n, g)$  becomes a singular manifold. Currently, the validity of the interpretation of  $D_\phi$  as a canonical divergence in singular manifolds is not clear. The generalisation of canonical divergences and of Theorem 3.1 to singular manifolds is currently under investigation.

### 3.3.3 Models involving constraints on the parameter space

Finally, in cases where  $x|y$  is constrained to a convex set  $\mathcal{S} \subset \mathbb{R}^n$  only the first and the third results of Theorem 3.1 hold. Proceeding similarly to the proof of Theorem 3.1, it is easy to show that  $D_\phi$  is the canonical divergence of the manifold  $(\mathcal{S}, g)$ , and that the Bayesian estimator related to the dual divergence is  $\hat{x}_{MMSE} = \operatorname{argmin}_{u \in \mathcal{S}} \mathbb{E}_{x|y}[D_\phi^*(u, x)]$ . However, the Bayesian estimator that minimises the canonical divergence is now a shifted MAP estimator

$$\hat{x}_{D_\phi} = \operatorname{argmin}_{u \in \mathcal{S}} D_\phi(u, \hat{x}_{MAP}) + u^\top \mathbb{E}_{x|y}[\nabla \phi(x)],$$

where generally  $\mathbb{E}_{x|y}[\nabla \phi(x)] \neq 0$  (see the appendix for details). Therefore, for these models  $\hat{x}_{MAP}$  is potentially not a proper Bayesian estimator in the decision-theoretic sense. Nevertheless, the high dimensional guarantees of Propositions 3.1 and 3.2 still hold for  $\hat{x}_{MAP}$ , providing partial theoretical justification for using this estimator (also observe that  $\hat{x}_{MAP}$  is in the neighbourhood of  $\hat{x}_{MMSE}$  in the sense of Proposition 3.2).

## 4 Conclusion

MAP estimation is one of the the most successful Bayesian estimation methodologies in modern data science, with a track record of accurate results across a wide range of challenging applications involving very high dimensionality. Our aim here has been to contribute to the theoretical understanding of this widely used methodology, particularly by placing it in the Bayesian decision theory framework that underpins the core Bayesian inference methodologies.

In order to analyse MAP estimators we have adopted an entirely new approach: we allowed the model to self-specify the loss function, or equivalently the Bayesian estimator, that is used to summarise the information that the model represents. This was achieved by using the connections between model log-concavity, Riemannian geometry, and divergence functions. We first established that if  $p(x|y)$  is strongly log-concave, continuous, and  $\mathcal{C}^3$  on  $\mathbb{R}^n$ , then  $\phi(x) = -\log p(x|y)$  induces a dually-flat Riemannian structure on the parameter space, where the canonical divergence is the Bregman divergence associated with  $\phi$ , and where the MAP estimator is the unique Bayesian estimator w.r.t. to this loss function. We also established

that the MMSE estimator is the Bayesian estimator w.r.t. the dual canonical loss, and that both estimators enjoy favourable stability properties in high dimensions. We then examined the effect of relaxing these assumptions to models that are not smooth, strictly but not strongly convex, or that involve constraints on the parameter space.

The theoretical results presented in this work provide several valuable new insights into MAP and MMSE estimation. In particular, both estimators stem from Bayesian decision theory and from the consideration of the geometry of the parameter space, and exhibit an interesting form duality. Also, the bounds on the expected estimation error and the large error probability for MAP estimators support the remarkable empirical performance observed in large scale settings, such as imaging and machine learning problems. The results also show that the predominant view of MAP estimators as hastily heuristic or approximate inferences, motivated only by computational efficiency, is fundamentally incorrect (though the fact that MAP estimators are available as solutions to convex minimisation problems is a fundamental practical advantage). We hope that these new theoretical results will promote a wider adoption of this powerful Bayesian point estimation methodology across all domains of statistical data science.

## 5 Acknowledgements

Part of this work was conducted when the author held a Marie Curie Intra-European Research Fellowship for Career Development at the School of Mathematics of the University of Bristol. He is grateful to Yoann Altmann, Gavin Gibson, Peter Green, Abderrahim Halimi, Bernd Schroers, Jonty Rougier, and Ben Powell for useful discussion.

## References

- [1] S.-I. Amari and H. Nagaoka. *Methods of Information Geometry (Translations of Mathematical Monographs)*. American Mathematical Society, Rhode Island, USA, 2007.
- [2] N. Ay and S.-I. Amari. A novel approach to canonical divergences within information geometry. *Entropy*, 17(12):7866, 2015.
- [3] Arindam Banerjee, Srujana Merugu, Inderjit S. Dhillon, and Joydeep Ghosh. Clustering with bregman divergences. *J. Mach. Learn. Res.*, 6:1705–1749, December 2005.
- [4] R. Bassett and J. Deride. Maximum a Posteriori Estimators as a Limit of Bayes Estimators. *ArXiv e-prints*, November 2016.
- [5] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer New York, 2011.

- [6] S. Bobkov and M. Madiman. The entropy per coordinate of a random vector is highly constrained under convexity conditions. *IEEE Trans. Info. Theory*, 57(8):4940–4954, Aug 2011.
- [7] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [8] Martin Burger and Felix Lucka. Maximum a posteriori estimates in linear inverse problems with log-concave priors are proper Bayes estimators. *Inverse Problems*, 30(11):114004, 2014.
- [9] E. J. Candes, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inf. Theory*, 52(2):489–509, Feb 2006.
- [10] Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Found. Comput. Math.*, 9(6):717, 2009.
- [11] V. Cevher, S. Becker, and M. Schmidt. Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics. *IEEE Signal Process. Mag.*, 31(5):32–43, Sept 2014.
- [12] Patrick L. Combettes and Jean-Christophe Pesquet. *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, chapter Proximal Splitting Methods in Signal Processing, pages 185–212. Springer New York, New York, NY, 2011.
- [13] A. P. Dawid. The geometry of proper scoring rules. *Ann. Inst. Stat. Math.*, 59(1):77–93, 2007.
- [14] P. J. Green, K. Łatuszyński, M. Pereyra, and C. P. Robert. Bayesian computation: a summary of the current state, and samples backwards and forwards. *Statistics and Computing*, 25(4):835–862, June 2015.
- [15] R. Gribonval. Should penalized least squares regression be interpreted as maximum a posteriori estimation? *IEEE Trans. Signal Process.*, 59(5):2405–2410, May 2011.
- [16] M. Pereyra. Maximum-a-posteriori estimation with bayesian confidence regions. *SIAM J. Imaging Sci.*, 2016. submitted.
- [17] M. Pereyra, P. Schniter, E. Chouzenoux, J.-C. Pesquet, J.-Y. Tournieret, A.O. Hero, and S. McLaughlin. A survey of stochastic simulation and optimization methods in signal processing. *IEEE. J. Selected Topics in Signal Process.*, 10(2):224–241, Mar. 2016.
- [18] C. P. Robert. *The Bayesian Choice (second edition)*. Springer Verlag, New-York, 2001.

## Proofs of Theorem 3.1 and Propositions 3.1 and 3.2

### Proof of Theorem 3.1

The first part of Theorem 3.1 follows directly from differential geometry and from the regularity properties of  $\phi$  (see [1] for an introduction to differential geometry). From differential geometry, under the conditions of Theorem 3.1,  $\phi$  induces a Riemannian metric on  $\mathbb{R}^n$  with coefficients

$$g_{i,j}(x) = \partial_i \partial_j \phi(x),$$

and where we note that  $g(x) \in \mathcal{S}_{++}^n$  from the strong convexity of  $\phi$ . Similarly, we have the affine connection coefficients

$$\Gamma_{i,j,k} = \partial_i \partial_j \partial_k \phi(x).$$

Moreover, because  $\phi$  is convex it endows  $(\mathbb{R}^n, g)$  with a dual affine coordinate system  $\eta$ , related to the primal coordinate system by the duality  $\eta_x = \nabla \phi(x)$  and  $x_\eta = \nabla \phi^*(\eta)$ , where  $\phi^*(\eta) = \max_{x \in \mathbb{R}^n} x^\top \eta - \phi(x)$  is the convex conjugate of  $\phi$  [1, Ch. 3]. As a result we have a dual Riemannian metric  $g^*$  w.r.t.  $\eta$ , with coefficients given by

$$g_{i,j}^*(\eta) = \partial_i \partial_j \phi^*(\eta),$$

and a dual affine connection  $\Gamma^*$  with coefficients given by

$$\Gamma_{i,j,k}^*(\eta) = \partial_i \partial_j \partial_k \phi^*(\eta).$$

Finally, it is easy to check that  $x$  and  $\eta$  are mutually dual w.r.t.  $g$ . That is, for all  $x \in \mathbb{R}^n$

$$g^*(\eta_x) = g(x)^{-1}$$

which implies that  $(\mathbb{R}^n, g, \Gamma, \Gamma^*)$  is a dually-flat Riemannian manifold [1, Ch. 3].

From [2], in such manifolds the  $\Gamma$ -geodesic connecting  $u \rightarrow x$  in (2) is given by  $\gamma_t = u + t(x - u)$ , and  $\dot{\gamma}_t = x - u$ . The proof is then concluded by integration by parts of (2) to obtain the Bregman divergence  $D_\phi(u, x) = \phi(u) - \phi^*(\eta_x) - \eta_x^\top u$ , which also admits the more familiar expression  $D_\phi(u, x) = \phi(u) - \phi(x) - \nabla \phi(x)(u - x)$ .

To prove the second part of Theorem 3.1 we use the linearity property of the expectation operator to express the definition  $\hat{x}_{D_\phi} = \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(u, x)]$  as follows

$$\begin{aligned} \hat{x}_{D_\phi} &= \operatorname{argmin}_{u \in \mathbb{R}^n} \phi(u) + \mathbb{E}_{x|y}[\phi(x)] - u^\top \mathbb{E}_{x|y}[\nabla \phi(x)] - x^\top \mathbb{E}_{x|y}[\nabla \phi(x)], \\ &= \operatorname{argmin}_{u \in \mathbb{R}^n} \phi(u) - u^\top \mathbb{E}_{x|y}[\nabla \phi(x)]. \end{aligned}$$

In a manner akin to [8], the proof is concluded by using the divergence theorem, together with the fact that  $p(x|y)$  is continuous and vanishes at least exponentially as  $\|x\| \rightarrow 0$ , to show that  $\mathbb{E}_{x|y}[\nabla\phi(x)] = \int_{\mathbb{R}^n} \nabla p(x|y) dx = 0$ . Hence,

$$\begin{aligned}\hat{x}_{D_\phi} &= \operatorname{argmin}_{u \in \mathbb{R}^n} \phi(u), \\ &= \hat{x}_{MAP}.\end{aligned}$$

Note that in the case where  $p(x|y)$  involves hard constraints on the parameter space then generally  $\mathbb{E}_{x|y}[\nabla\phi(x)] \neq 0$ , and we have  $\hat{x}_{D_\phi} = \operatorname{argmin}_{u \in \mathbb{R}^n} D_\phi(u, \hat{x}_{MAP}) - u^\top \mathbb{E}_{x|y}[\nabla\phi(x)]$  generally different from  $\hat{x}_{MAP}$ .

Finally, the proof of the third part of Theorem 3.1 follows directly from [3, Proposition 1], which for completeness we detail below

$$\begin{aligned}\hat{x}_{D_\phi^*} &= \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi^*(u, x)], \\ &= \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(x, u)], \\ &= \operatorname{argmin}_{u \in \mathbb{R}^n} \mathbb{E}_{x|y}[D_\phi(x, u)] - \mathbb{E}_{x|y}[D_\phi(x, \hat{x}_{MMSE})], \\ &= \operatorname{argmin}_{u \in \mathbb{R}^n} \phi(\hat{x}_{MMSE}) - \phi(u) - (\hat{x}_{MMSE} - u)^\top \nabla\phi(u), \\ &= \operatorname{argmin}_{u \in \mathbb{R}^n} D_\phi(\hat{x}_{MMSE}, u), \\ &= \hat{x}_{MMSE}.\end{aligned}$$

### Proof of Proposition 3.1

Assume that  $\phi(x) = -\log p(x|y)$  is convex on  $\mathbb{R}^n$  and define the generalised divergence

$$D_{\phi,q}^*(u, x) = \phi(x) - \phi(u) - q^\top (x - u)$$

where  $q \in \mathbb{R}^n$  belongs to the subdifferential set of  $\phi$  at  $u$ , which collapses to  $q = \nabla\phi(u)$  when  $\phi$  is differentiable at  $u$  [5]. Then, setting  $q = 0$  from the optimality condition of  $\hat{x}_{MAP}$

$$\mathbb{E}_{x|y} \left[ \frac{D_{\phi,0}^*(\hat{x}_{MAP}, x)}{n} \right] = \mathbb{E}_{x|y} \left[ \frac{\phi(x)}{n} \right] - \frac{\phi(\hat{x}_{MAP})}{n}.$$

Noting that  $\mathbb{E}_{x|y} \left[ \frac{\phi(x)}{n} \right]$  is the entropy rate of  $x|y$ , we use Proposition I.2 of [6] and obtain

$$\mathbb{E}_{x|y} \left[ \frac{D_{\phi,0}^*(\hat{x}_{MAP}, x)}{n} \right] \leq 1.$$

Finally, it follows from the proof of Theorem 3.1 that when  $\phi \in \mathcal{C}^1$ ,  $\hat{x}_{MMSE}$  minimises the posterior expectation of  $D_{\phi,q}^*(\hat{x}_{MMSE}, x)$  with  $q = \nabla\phi(\hat{x}_{MMSE})$ ,

which in turn implies that

$$\mathbb{E}_{x|y} \left[ \frac{D_{\phi,q}^*(\hat{x}_{MMSE}, x)}{n} \right] \leq \mathbb{E}_{x|y} \left[ \frac{D_{\phi,0}^*(\hat{x}_{MAP}, x)}{n} \right] \leq 1.$$

### **Proof of Proposition 3.2**

The proof follows directly from using Theorem 1 of [16] to express the set  $\{x : D_{\phi,0}^*(\hat{x}_{MAP}, x)/n < 1 + \epsilon\}$  as a posterior high-posterior-density credible region of level  $(1 - \alpha_\epsilon) \geq 1 - 3 \exp\{-n\epsilon^2/16\}$ .